CaFA: Cost-aware, Feasible Attacks With Database **Constraints Against Neural Tabular Classifiers** 

Matan Ben-Tov, Daniel Deutch, Nave Frost, Mahmood Sharif 😭 😥





## Motivation

ML on **tabular** data is ubiquitous in many critical applications.



**Caveat:** ML algorithms are susceptible to **evasion attacks**.

These attacks are thoroughly explored in vision, and can even be **realized** (as physical artifacts) generically.





## Can we generically evaluate robustness against realizable evasion attacks in the tabular domain?

Naïve attempt, apply the same attack: Main challenge – problem-feature space gap (heterogenous mapping that is neither differentiable nor invertible) [1]

Problem Space Feature Space PctNull | PctExt isMailto Hyperlinks Hyperlinks Phishing website

## **Background and Prior Work**

**Existing attacks** (mainly) lack in realizability or genericness:

- 1. Problem-Space Attacks Manipulate problem-space instances *directly* via allowed transformations. E.g., Lucas et al. 2021, Eykholt et al. 2023
- 2. Feature-Space Attacks Manipulate feature-space instances + accounting for realizability and imperceptibility.

E.g., Simonetto et al. 2022, Mathov et al. 2022, Sheatsley et al. 2021, Kireev et al. 2023

For **realizability**, we consider **data-integrity constraints types**:

Ň

$x' \not\models \forall x: x. is Mailto \in \{0,1\}$					Structure Constraints
	Structure can express			express	Involve structural properties
		PctNull PctExt		Feature's type, permissible range, etc.	
	ISMailto	Hyperlinks	Hyperlinks		Relational Constraints



